

Modelling representative population mobility for COVID-19 spatial transmission in South Africa

A. Potgieter¹, I. N. Fabris-Rotelli^{1,*}, Z. Kimmie², N. Dudeni-Tihone³,
J. Holloway³, C. Janse van Rensburg⁴, R. Thiede¹, P. Debba^{3,5},
R. Manjoo-Docrat⁵, N. Abdelatif⁴, and S. Makhanya⁶

¹University of Pretoria, Department of Statistics, South Africa

²Foundation of Human Rights, South Africa

³Council for Scientific and Industrial Research, South Africa

⁴Biostatistics Research Unit, South African Medical Research Council, South Africa

⁵Department of Statistics and Actuarial Science, University of Witwatersrand, South Africa

⁶IBM Research, Johannesburg 2001, South Africa & College of Graduate Studies, Unisa

Correspondence*:

Inger Fabris-Rotelli

inger.fabris-rotelli@up.ac.za

2 ABSTRACT

3 The COVID-19 pandemic starting in the first half of 2020 has changed the lives of everyone
4 across the world. Reduced mobility was essential due to it being the largest impact possible
5 against the spread of the little understood SARS-CoV-2 virus. To understand the spread, a
6 comprehension of human mobility patterns is needed. The use of mobility data in modelling is
7 thus essential to capture the intrinsic spread through the population. It is necessary to determine
8 to what extent mobility data sources convey the same message of mobility within a region. This
9 paper compares different mobility data sources by constructing spatial weight matrices at a
10 variety of spatial resolutions and further compares the results through hierarchical clustering. We
11 consider four methods for constructing spatial weight matrices representing mobility between
12 spatial units, taking into account distance between spatial units as well as spatial covariates.
13 This provides insight for the user into which data provides what type of information and in what
14 situations a particular data source is most useful.

15 **Keywords:** COVID-19, spatial, mobility, spatial weight matrices, principal component analysis, hierarchical clustering

1 INTRODUCTION

16 The COVID-19 pandemic starting in the first half of 2020 has changed the lives of everyone across the
17 world. From working from home at all hours, using less public and personal transport, home-schooling
18 under lock down, to economic strife and anxiety; predicting such changes would have been near impossible
19 a priori. By far the largest impact, aside from the economic troubles many find themselves in, is reduced
20 mobility. Daily commuting has been much reduced due to various lockdown measures internationally. In

21 addition, international flights and cross border travel was restricted for significant periods of time, even
22 between regions in some countries.

23 Reduced mobility was essential, however, due to it being the largest impact possible against the spread of
24 the little understood SARS-CoV-2 virus. Social distancing and stay at home instructions were understood
25 and implemented internationally. These instructions were seen as the best protection for the individual,
26 as well as being the means to prevent overload on the hospital systems, which would otherwise result in
27 inflated death rates. These protection mechanisms are formed on an understanding of the basic nature of the
28 spatial spread of the virus. A virus spreads via a host, whom it relies on to move amongst other susceptibles.
29 The more movement and interaction performed by the host, the more the virus is able to spread. It is thus
30 imperative to incorporate a spatial element when modelling the spread of the COVID-19 pandemic. Herein,
31 we focus on modelling the mobility spatially.

32 Quantifying mobility patterns of people facilitates a more accurate understanding of the spread of the
33 disease. An individual's ability to physically "lock down" and stay at home was affected by economic
34 inequality, as shown in a US study [?]. In South Africa, this economic inequality is extreme, with the
35 World Bank recognising South Africa, in 2019, as having the worst inequality in the world¹.

36 While the strict lockdown introduced by the South African government from 27 March 2020 delayed the
37 first wave, the mobility was by no means completely reduced due to many living day-to-day for food. Food
38 parcel queues from food donations were a large focus during the first half of the pandemic in South Africa,
39 as the risk of contracting COVID-19 was overridden by the need for food. Such queues, and the use of
40 public transport during these times, heightened the transmission risk of COVID-19 in South Africa, even
41 while lockdown rules were in place. A full lockdown was therefore not possible, and spatial interaction
42 continued between individuals from different regions across South Africa. Modelling regions in isolation
43 will therefore not capture the influence of this mobility on the spread of COVID-19 in South Africa.

44 The use of mobility data in modelling COVID-19 is thus essential to capture the intrinsic spread through
45 the population. A common source is mobile phone location data, which has been utilized previously for
46 epidemiological modelling [? ? ? ? ?]. However, this data is difficult to obtain due to increasing privacy
47 concerns worldwide. In addition, there are often a number of network providers in a region, each with
48 certain market share. Without data access from all, or at least, the largest providers, representativeness and
49 mobile phone penetration will be limited and should be used with caution. Other sources of mobility data
50 are published by Facebook and Google. The spatial resolution of these is lower, however. In this paper we
51 focus on mobile phone and Facebook mobility data, which has higher spatial resolution than the Google
52 alternative.

53 It is necessary to determine to what extent different sources of mobility data, at differing spatial resolutions,
54 convey the same message of mobility within a region. In this paper we demonstrate, through the use of
55 principal component analysis as well as hierarchical clustering, how different sources of spatial mobility
56 data at various resolutions can lead to different conclusions with regards to spatial unit connectivity.
57 Spatial connectivity is an essential first step in spatial modeling, providing a quantification of the spatial
58 dependency between spatial units. Herein, we compare the calculation of a number of spatial weight
59 matrices in quantifying how spatial units relate. We also discuss the advantages of different sources and
60 how they can be harnessed when modelling the spread of a virus. We do this by using principal component
61 analysis in order to condense the information that can be gained from a spatial weight matrix and then
62 using hierarchical clustering to identify the strongest spatial associations and to essentially put on display

¹ <https://povertydata.worldbank.org/Poverty/Home> (Accessed May 2021)

63 what type of relationships the spatial weight matrix is identifying. This is to the best of our knowledge the
64 first time this exact combination has been used for this purpose.

65 The mobility data available for South Africa is presented in Section 2. The methodology for constructing
66 connectivity matrices is developed in Section 3, and the results are presented in Section 4. Section 5
67 provides a discussion and Section 6 concludes.

2 DATA

68 Available mobility data is at different resolutions. For the case of South Africa, the administrative divisions
69 of the country are summarised in Table 1. In order of increasing spatial resolution these are country,
70 province, district municipality, local municipality, and ward, labelled as administrative levels 0 through 4
71 respectively. To facilitate the comparison of different sources of spatial information, it is first necessary to
72 aggregate the data from each source to the same spatial resolution. Increasing the resolution of spatial data
73 can be achieved through methods such as small area estimation or spatial micro-simulation (see e.g. [? ?]).
74 These methods are somewhat involved and require the use of auxiliary information or assumptions that are
75 unlikely to be true. In this paper we investigate aggregating down to the lowest spatial resolution used by
76 our data sources. While this is relatively straightforward to accomplish, it potentially results in the loss of
77 information.

Table 1. South Africa's administrative boundaries

Administrative level	Spatial unit name	Number of spatial units
0	Country	1
1	Province	9
2	District municipality	52
3	Local municipality	213
4	Ward	4392

78 Mobility data are used to understand various issues ranging from epidemic modelling, transport plan-
79 ning and management, communication network improvement and urban planning [? ?]. Asgari et al.
80 [?] indicates that mobility goes far beyond mere geographical movement of humans, but provides a
81 comprehensive perspective on human interactions that could be considered from spatial, temporal, and
82 contextual aspects. Human mobility is one of the aspects of mobility that gained attention from the global
83 spread of infectious diseases as with the recent COVID-19 pandemic. A variety of technologies including
84 navigation sensors, wireless technologies, and cellular communication technologies are used to position
85 humans in space [?]. A study by Zhou et al. [?] provides a comprehensive overview of the different types
86 of human mobility patterns data. These include those data types that capture both the wider (city-wide)
87 and minute (building-wide or large structure) human movements, for example, cellular services records
88 (CSRs), surrounding WiFi access point records (SWAPRs), Global Positioning System locations (GPSLs),
89 geotagged social media (GTSM), public transport smart card records (PTSCRs), bluetooth detection records
90 (BDRs), and WiFi probe request records (WFPRs). The analysis methods range from data visualisation to
91 statistical analysis methods (classification and clustering), heuristic logic, graph theory and optimization

92 techniques.

93 **2.1 South Africa's lockdown levels**

94 To quell the spread and impact of the COVID-19 pandemic, the South African government instigated one
 95 of the strictest lockdowns in the world. This particular lockdown strategy is structured around different
 96 "levels" of lockdown, each of which brings different restrictions (with level 5 being the highest and placing
 97 restrictions on nearly all forms of travel to all citizens except for those classified as essential workers). The
 98 various levels as well as the dates for which they were active are given in Table 2. Note that for this paper
 99 we only consider the lockdown until the end of Level 3 due to data availability only over this period.

Table 2. South Africa's lockdown levels and dates

Level	Date	Restrictions
Business as usual	1 March 2020 - 26 March 2020	No restrictions
Level 5	27 March 2020 - 30 April 2020	Essential services only otherwise all confined to place of residence. No inter-provincial movement, except for transportation of goods and exceptional circumstances e.g. funerals. Public and private transport restricted to certain times of the day, with limitations on vehicle capacity
Level 4	1 May 2020 - 31 May 2020	More sectors permitted with restrictions, including mining, and partial e-Commerce allowed. Public places (such as religious, cultural, recreational facilities) and the tourism sector remain closed and gatherings prohibited. All confined to place of residence from 8pm-5am. No local (between metropolitan areas or districts) or inter-provincial movement of people, except for permitted reasons e.g. returning for alert level 4 operations. All borders remain closed except for designated ports of entry for restricted home affairs operations and for the transportation of fuel, cargo and goods. Public and private transport may operate at all times of the day, with limitations on vehicle capacity
Level 3	1 June 2020 - 17 August 2020	More sectors permitted including take away restaurants, e-commerce and delivery services and global business services. Public places and tourism opened and gatherings and sporting activities permitted but all subject to restrictions. All confined to place of residence from 11pm-4am. No inter-provincial movement of people, except for transportation of goods, exceptional circumstances and other permitted reasons. Public and private transport may operate at all times of the day, with limitations on vehicle capacity

100 As non-pharmaceutical interventions (such as the lockdown) are eased the population is allowed to

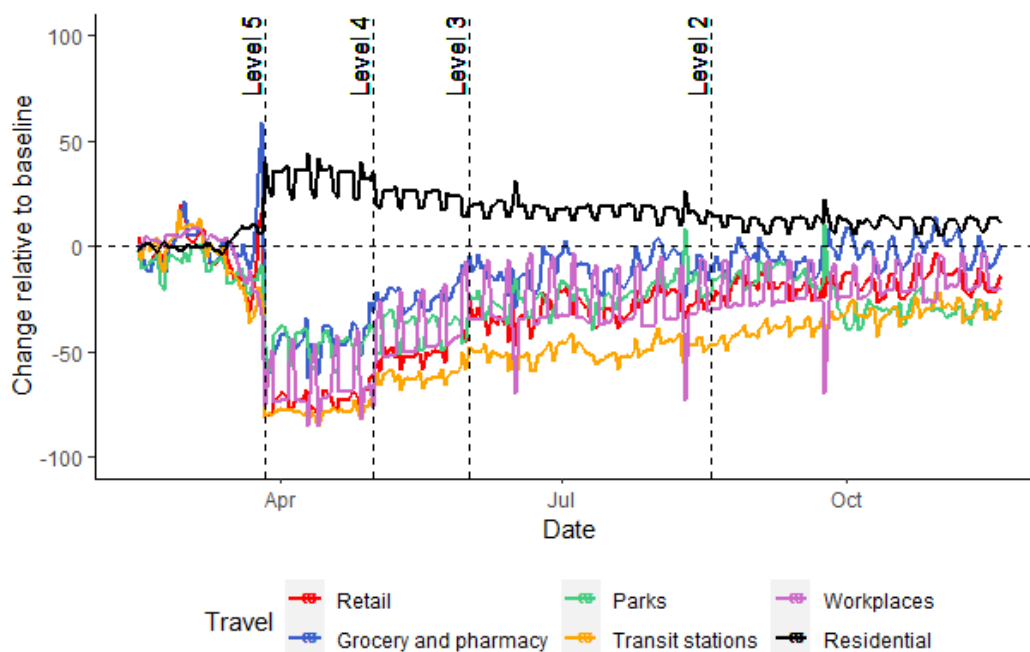


Figure 1. Google mobility report data for 15 February 2020 – 20 November 2020 (transitions to different levels of lockdown indicated by vertical reference lines)

101 become more mobile. Naturally this will have an impact on the transmission rate of the virus and thus this
 102 temporal element must be included in some manner. In this paper we split the data temporally on the date
 103 ranges given in Table 2 up to level 4 and set up a spatial weight matrix for each level of lockdown to study
 104 how mobility patterns changed.

105 Two mobility data types were available for this research. The first is freely available data shared by
 106 Facebook, and the second is mobility data made available by a South African cellular provider for the
 107 context of the COVID-19 response in 2020. In Figure 7 we provide the Google mobility data at country
 108 level. We do not use this data in this research as it is only available at administrative level 1, representing
 109 low spatial resolution. It is however useful for context providing mobility levels in each various industry
 110 sectors. Mobility for residential travel (i.e., individuals remaining at their place of residence) is the only
 111 type of travel that saw an increase after the country transitioned into level 5. Grocery and pharmacy travel
 112 saw an initial spike shortly before the country went into level 5 (possibly attributed to panic-buying). After
 113 transitioning to level 5 we see a drastic decrease in all types of travel, with residential travel showcases a
 114 slightly downward trend while all other forms of travel have an upward trend. Grocery and pharmacy travel
 115 is the quickest to recover to pre-COVID levels while travel to parks and travel stations is the slowest to
 116 recover (most likely due to this being for leisure). By the end of the year residential travel is still higher
 117 than before any lockdown interventions. Table 3 provides the average changes over each level as well.

118 2.2 Facebook Data for Good

119 Multiple geographically indexed datasets have been made freely available for use by Facebook through
 120 their “Facebook data for good” initiative. These datasets serve to aid researchers and policymakers in
 121 understanding the spread of COVID-19².

122 This paper utilises one of these available datasets, namely the “Movement range maps” dataset. The data

² <https://dataforgood.fb.com/> (Accessed May 2021)

Level	Date	Retail	Grocery and pharmacy	Parks	Transit stations	Workplaces	Residential
BAU	2 Feb - 26 Mar	-3.49	1.68	-9.39	-5	-0.88	1.71
Level 5	27 Mar - 30 Apr	-73.06	-46.09	-46.86	-78.49	-65.89	33
Level 4	1 May - 31 May	-50.39	-23.45	-39.39	-61.71	-40.58	23.35
Level 3	1 Jun - 17 Aug	-29.53	-10.71	-23.17	-49.72	-28.1	17.17
Level 2	18 Aug onwards	-17.76	-3.34	-23.29	-34.65	-19.78	11.35

Table 3. Average changes in population mobility over lockdown levels using the Google mobility data during 2020

123 indicates the change in mobility, $F_i^{(t)} \in (-1, 1)$ (which can be interpreted as a percentage $(-100, 100)$),
 124 for a spatial unit i on a given day t over the period 1 March 2020 – 28 February 2021 relative to a one-week
 125 baseline calculated in February 2020. The daily values for each district municipality were calculated by
 126 determining the number of so-called “Bing tiles”³ that each inhabitant visited on a given day (place of
 127 residence being determined by the location where users most often spend their nights). A Bing tile is the
 128 term used by Microsoft for a spatial polygon. After incorporating some degree of noise, the average number
 129 of tiles visited by the inhabitants was determined and expressed relative to the baseline. The full description
 130 of how these values were calculated is available in the Appendix. The spatial resolution for units of this
 131 data are district municipalities, namely at administrative level 2.

132 Figure 2 shows the aggregated data for district municipalities, with the average across the district munici-
 133 palities shown in red. The figure demonstrates that the average mobility nationally dropped significantly in
 134 late March. This corresponds to when South Africa entered its first hard lockdown on the 27th of March
 135 2020 (see Table 2). The hard lockdown imposed severe restrictions on travel and constituted a strict stay
 136 at home directive. Only essential workers were allowed to leave their homes. Furthermore, the average
 137 change in mobility is primarily negative over the entire study period, indicating that mobility patterns
 138 remain more constrained than before the hard lockdown. The first COVID-19 case was discovered on 5
 139 March 2020 and the lockdown announcement was made a week later on 15 March. This could explain the
 140 drop in mobility already seen from early March.

141 Notable benefits of using this data are that the data is freely available and could potentially act as a
 142 very representative proxy for human mobility, as Facebook services are not constrained to specific mobile
 143 network providers. In addition, all the cellular network providers in South Africa provide a free version
 144 of Facebook called Facebook Zero. Even though it is known that not all South Africans have a Facebook
 145 account, the Facebook mobility data may provide an acceptable level of representativeness for mobility
 146 within the country since the population of South Africa is considered significantly young⁴. It is also clear
 147 that a large amount of the original data was censored in order to preserve user privacy and thus the data is at
 148 a sparse level of spatial resolution (administrative level 2). The data is also not specific with regards to the
 149 direction of spatial mobility. Daily observations only indicate whether individuals were more or less mobile

³ <https://docs.microsoft.com/en-us/bingmaps/articles/bing-maps-tile-system> (Accessed May 2021)

⁴ Mid-2021 Statistics South Africa Population Report <http://www.statssa.gov.za/publications/P0302/P03022021.pdf> (Accessed August 2021)

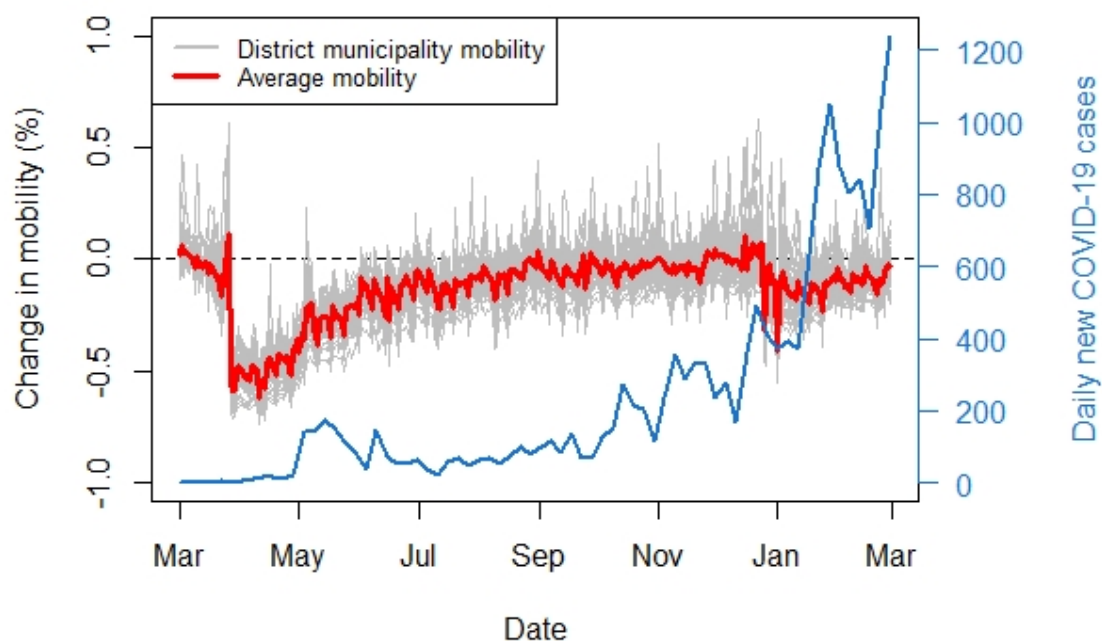


Figure 2. “Facebook for good” movement range maps data (1 March 2020 – 28 February 2021) relative to a baseline calculated in a week of February 2020

150 in a district municipality and do not indicate the spatial units towards which this mobility was directed.

151 **2.3 Mobile network data**

152 The growing popularity and widespread use of mobile devices has led to massive amounts of data being
 153 produced at any given point in time all around the world. Mobile phone data can be collected either
 154 passively by mobile services providers or through the use of mobile applications. The ease with which
 155 such large quantities of data can be gathered makes cellular data attractive for researchers. Mobile devices
 156 operate by sending and receiving information from cellphone towers. When interacting with a cellphone
 157 tower we say that a phone has “pinged” off a cellphone tower. A mobile device may ping off a cellphone
 158 tower by sending or receiving any kind of information, be it a phone call, text message or application
 159 notification. The mobile network data obtained for this research is obtained using the number of users
 160 whose mobile devices pinged off a cellphone tower within one ward (administrative level 4) on a given day
 161 and then later that day pinged off a cellphone tower in a different ward.

162 Mobile phone data has been used numerous times in the field of spatial epidemiology to model the spread
 163 of various diseases, including cholera [? ?], dengue [? ?] and malaria [? ?]. Following the outbreak
 164 of the COVID-19 pandemic, the governments of various countries across the world began collecting
 165 cellular device user data in an attempt to aid the conception and implementation of non-pharmaceutical
 166 interventions [? ? ? ?]. This data has since been used by researchers to clearly establish a correlation
 167 between population mobility and COVID-19 case numbers [? ? ? ?].

168 Limitations of mobile phone data exist. First and foremost of these is the issue of user privacy. Mobile
 169 phone data could potentially be misused to identify specific individuals and thus cellular providers are

170 often hesitant to provide researchers with such data [? ?]. Such data is often aggregated to a low spatial
 171 resolution to prevent this as well as reduce noise, but this comes at the cost of some data specificity. Another
 172 potential drawback of mobile phone data is high computational cost. For high mobile phone penetration
 173 rates, mobile phone data may consist of a number of entries in the order of billions. The computational cost
 174 of processing such datasets is prohibitive, potentially preventing analysis.

175 For this paper, anonymised mobile phone data was obtained from a local mobile network provider. In
 176 South Africa, the mobile phone penetration level is estimated to be as high as 95%⁵. The provider utilised
 177 in this paper is one of the largest providers in the country, with an estimated market share of 42%.

The data provides the number of mobile phone users $m_{ij}^{(t)}$ that travelled to ward j from ward i on day t for the period 2 March - 12 May 2020. The data is at administrative level 4, which is the highest spatial resolution reasonably possible while preserving some level of privacy of exact user location. To compare insights gained from this data and the Facebook dataset in Section 2.2, it would first be necessary to aggregate the mobile phone data to the same spatial resolution which is administrative level 2. In South Africa, each ward has a unique 8-digit ID code. The first three digits of this code indicates the district municipality that the ward is a part of. For example, the ward ID 9344007 indicates that the ward is part of the district municipality with code 934. In order to aggregate the data to district municipality level, one could replace the ward IDs of the observations with their district municipality codes (i.e. only the first 3 digits), whereupon rows with identical origin and destination codes would be discarded. The mobile phone data at administrative level 2 is thus given by

$$M_{I,J}^{(t)} = \sum_{i \in I, j \in J} m_{ij}^{(t)},$$

178 where I and J are district municipalities and i and j are wards as previously indicated. Transitions contained
 179 within a single district municipality are thus discarded. Analysis revealed that this caused an average of
 180 26% of daily observations to be discarded. The retained data is displayed in Figure 3. The representation
 181 differs to that of Figure 1 as the data provides transitions between regions in this case. We once again notice
 182 a sharp decline in population mobility in late March.

183 The population of South Africa (mid-2021) is approximately 60.14 million⁶, and yet the highest total
 184 number of inter-district municipality transitions on any given day was approximately 10 million (seen in
 185 Figure 3). It should be noted that the same individual can be responsible for multiple transitions and that
 186 some individuals could potentially possess multiple cellular devices. Literature does exist on the use of
 187 mobile phone data to estimate population numbers, see e.g. [?]. Doing so is not within the scope of the
 188 research presented here but would be of value in testing mobile phone representability.

189 Despite the quality of available hardware⁷, this process proved highly computationally expensive due to
 190 the number of comparisons that need to be run on billions of lines of data in order to create a spatial weight
 191 matrix for each day in the time period.

3 METHODOLOGY

⁵ See <https://www.geopoll.com/blog/mobile-penetration-south-africa/> and <https://www.icasa.org.za/uploads/files/State-of-the-ICT-Sector-Report-March-2020.pdf> (Accessed May 2021)

⁶ Mid-2021 Statistics South Africa Population Report <http://www.statssa.gov.za/publications/P0302/P03022021.pdf> (Accessed August 2021)

⁷ All analysis presented here was performed on a desktop computer running Intel Core i7 with a clock speed of 3.40GHz, a 64-bit operating system and 64 GB of installed memory.

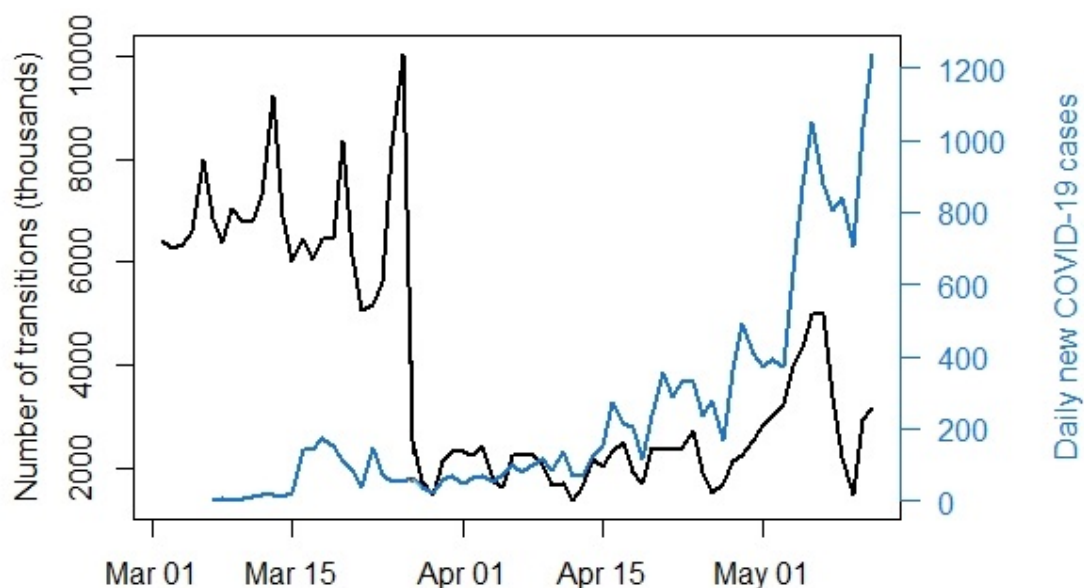


Figure 3. Number of individual transitions between wards using the available mobile phone data (2 March 2020 – 12 May 2020)

192 3.1 Literature review

193 When a particular phenomenon exhibits evidence of spatial dependence, this dependency must be taken
 194 into account when modelling to minimise the risk of producing biased results [? ?]. In the case of an
 195 infectious disease that is spread through physical contact and near proximity, it is clear that locations
 196 that are situated closer together (or rather the inhabitants of these locations) will play a larger role in
 197 determining their respective infection rates than locations that are farther apart. To incorporate this fact,
 198 spatial models allow spatial units to be more strongly (or weakly) correlated with one another based on
 199 some select criteria that is deemed suitable for the phenomenon being modelled. This is achieved through
 200 the use of a spatial weight matrix (sometimes called a “spatial mobility matrix”) usually denoted by \mathbf{W} [?
 201 ? ? ? ? ?].

202 **Definition 1** (Spatial weight matrix). Let $S = \{1, 2, \dots, n\}$ be a set of spatial units. A spatial weight
 203 matrix [? ? ? ? ?] is an $n \times n$ matrix $\mathbf{W} = [w_{ij}]$ satisfying $w_{ij} \geq 0$ and $\sum_{j=1}^n w_{ij} = 1 \quad \forall i \in S$.

204 This matrix is formally defined as an expression of spatial dependency between spatial units [? ? ? ? ?].
 205 Simply put, the spatial weight matrix is constructed in such a way so that entry w_{ij} quantifies the amount
 206 of spatial influence that spatial unit i exerts on spatial unit j [? ? ? ? ?].

207 Such matrices are frequently restricted to being symmetrical to simplify estimation. However, symmetry
 208 is not required and can result in a less realistic representation of spatial dependency [? ? ? ? ?]. Another
 209 common convention is that $w_{ii} = 0$ for all i to exclude the possibility of so-called “self-influence” [? ? ? ? ?].
 210 Non-zero diagonal entries can however be included and are interpreted as quantifying the resistance that

211 each spatial unit has against influence from the other spatial units [? ?]. Performing row-standardisation
212 on the matrix allows the connectivity of different spatial units to be compared [? ?].

213 Spatial weight matrices are most commonly used in the fields of econometrics and spatial statistics [?].
214 Recently however, they have become popular in the field of spatial epidemiology and have been used to
215 model various diseases including dengue, malaria, foot and mouth disease [? ? ? ?] and most recently
216 COVID-19 [? ?]. There are relatively few established guidelines with regards to constructing a spatial
217 weight matrix [? ? ? ?], however, the construction of these matrices has seen some advancement, with
218 greater emphasis being placed on creating matrices that offer an accurate representation of human mobility.
219 Simpler models rely on measures such as distance, contiguity or adjacency [? ? ? ? ? ? ? ?] while
220 more complex ones are able to use mobile phone data [?] and geostatistical information [? ?]. Accurately
221 specifying these matrices is a non-trivial problem, as discussed in [?]. Most recently, Ejigu et al. proposed
222 a methodology through which both distance and covariate information can be utilized [?].

223 Given the importance of correctly specifying the spatial weight matrix, and the fact that there are
224 often multiple sources of spatial data available on hand, it becomes necessary to develop some means of
225 comparing spatial weight matrices. Specifically, it is necessary to compare the insights that can be derived
226 from different spatial weight matrix definitions. In recent years this comparison has been achieved either
227 through the use of measures of spatial autocorrelation, such as Moran's I [?], or through more specialised
228 methods local to the field of spatial statistics [? ?]. In this paper, we adapt an idea initially presented
229 by Garrison and Marble [?], whereby principal component analysis is used to reduce the dimensionality
230 of candidate spatial weight matrices. We then introduce the use of hierarchical clustering to derive a
231 clustering solution for the spatial unit principal scores. This allows for a more informative comparison of
232 the information provided by these connectivity matrices, as opposed to simply comparing their structure
233 visually.

234 3.2 Spatial weight matrices

235 Selecting an optimal spatial weight matrix is often reliant on the use of a priori information and experience.
236 In this paper the emphasis is on comparing the implications for different spatial weight matrices and
237 the varying types of spatial associations that they represent. We next discuss the spatial weight matrix
238 construction approaches used in this paper.

239 3.2.1 Method 1: Distance method

240 The exponential distance definition of a spatial mobility matrix is used frequently in studies involving
241 spatial correlation, and is a popular choice in spatial econometrics [? ? ? ?]. As previously mentioned
242 however, the concepts of distance, contiguity and adjacency do not necessarily offer the most accurate or
243 realistic representation of human mobility. In this paper we include this model in order to draw comparisons
244 between it and more data-driven models. The entries of the spatial weight matrix are given by

$$w_{ij} = \exp(-d_{ij}) \quad (1)$$

245 where d_{ij} is the Euclidean distance between the centroids of spatial units i and j . Diagonal entries are set
246 to 0 to remove the possibility of so-called "self-influence", and all rows are standardised to sum to 1 to
247 facilitate comparisons between different spatial units. Both of these restrictions were maintained for all
248 matrices in this paper. Under this model, spatial units are most strongly spatially correlated with the spatial
249 units that are closest to them geographically. No temporal component can be incorporated for this method.

250 3.2.2 Method 2: Mobile network method

251 The mobile network data indicates the number of individuals that travelled from spatial unit i to spatial
252 unit j on a given day t . These entries are used to construct the spatial weight matrix as follows,

$$w_{ij}^{(t)} = M_{ij}^{(t)}. \quad (2)$$

253 This model expresses spatial weights as a function of the amount of flux (both in and out) occurring at a
254 spatial location, and is sometimes referred to as a spatial interaction matrix [?]. Spatial units where more
255 (less) individuals travelled to other spatial units will thus have a larger (smaller) effect on other spatial
256 units.

257 3.2.3 Method 3: Weighted Facebook data method

258 In order to create a spatial mobility matrix using the Facebook data, we use the same approach of Ejigu
259 et al. [?]. This takes into account proximity as well as covariate information which is spatially dependent.
260 The entries of the the spatial weight matrix are given by

$$w_{ij}^{(t)} = \exp\left(-\left(\alpha \cdot |F_i^{(t)} - F_j^{(t)}| + (1 - \alpha) \cdot d_{ij}\right)\right) \quad (3)$$

261 where $F_i^{(t)}$ is the mobility of spatial unit i at time t , scaled by population size (the covariate information),
262 d_{ij} is the Euclidean distance between the centroids of spatial units i and j , and $\alpha \in (0, 1)$ is a control
263 parameter indicating the amount of weight that should be given to the covariate term. The control parameter
264 α was set to 0.6 in this paper to allow for the covariate data to play a slightly more prominent role in the
265 estimation process without disregarding the importance of distance. The parameter captures that we are
266 making an assumption that the Facebook data can be used to capture transitions between regions even
267 though it is isolated location data. The value of 0.6 gives the weighted calculation a slight nudge towards
268 the Facebook data. Note that if $\alpha = 0$ then the model simplifies to the exponential distance model in
269 equation (1).

270 The Facebook mobility data for each district municipality was scaled using population size in order to
271 account for the fact that increased mobility in a given district is more (less) influential to neighbouring
272 districts if the population size is large (small). This was also done in order to restore some of the variation
273 in the data that was likely lost when the data was censored to a lower spatial resolution.

274 3.2.4 Method 4: Scaled Facebook data method

275 An additional final spatial weight matrix was constructed based on further variation of the exponential
276 distance model. For this matrix, the rows of the exponential distance matrix are scaled using the (unscaled)
277 Facebook mobility data. For example, if the mobility within district municipality i was 20% lower than the
278 baseline, then the entire row i is multiplied by 0.8. Each entry in the exponential distance matrix is thus
279 scaled by some number in (0,2). The entries in the matrix are given by

$$w_{ij}^{(t)} = \left(1 + F_i^{(t)}\right) \cdot \exp(-d_{ij}). \quad (4)$$

280 This construction allows the exponential distance matrix to be scaled such that the spatial influence
281 of more (less) mobile district municipalities is increased (decreased). This also renders the exponential
282 distance matrix non-symmetric, which should offer a more realistic representation of spatial influence.

283 Methods 3 and 4 are a novel approach to constructing connectivity matrices from the Facebook mobility

284 data.

285 **3.3 Principal Component Analysis**

286 Principal component analysis (PCA) is a statistical technique that aims to derive a parsimonious repre-
287 sentation of a given dataset by deriving an orthogonal linear transformation of the data [?]. In standard
288 PCA, the only hyperparameter that needs to be selected is the number of principal components, which is
289 primarily dependent on the cumulative proportion of variance in the data that the user wishes to retain. For
290 this paper, the number of principal components was chosen such that 75% of the variation in the data was
291 maintained. The full discussion of PCA and its various extensions is left to the existing literature (see e.g.
292 [?]).

293 **3.4 Hierarchical clustering**

294 Hierarchical clustering is an unsupervised machine learning technique that allows the user to group
295 together data points in an attempt to uncover sets of observations that share similar characteristics [?].
296 This is achieved by procedurally grouping together those observations that are most similar to each other
297 based on some selected measure of dissimilarity, referred to as a “linkage” [?]. The number of retained
298 clusters can then be selected either using some measure of cluster (dis)similarity or a pre-selected value.
299 We use agglomerative clustering, which additionally requires the selection of a method through which the
300 dissimilarity of separate clusters is calculated. A full discussion on hierarchical clustering may be found in
301 [?].

302 Herein, we chose the number of clusters to be identical to the number of principal components. Complete
303 linkage was used to calculate the difference between clusters at each iteration. Single and average linkage
304 displayed a propensity for resulting in clusters that were very large. This was most likely due to the fact that
305 single linkage considers the minimum distance between clusters at each iteration, thus regarding clusters
306 as more similar in general. Complete linkage considers the maximum distance between clusters and thus
307 considers clusters to be more distinct. Average linkage is the average of these two extremes.

308 **4 RESULTS**

309 Figure 4 shows the 52 district municipalities of South Africa. The four largest cities in the country are
310 Tshwane, Johannesburg, Durban and Cape Town, situated in the City of Tshwane, City of Johannesburg,
311 eThekweni and City of Cape Town district municipalities respectively as indicated in colour in Figure 4.
312 These four cities are the focal point of economic activity and travel in the country, and it is thus logical that
313 they would play a substantially larger role in the transmission of the virus than other municipalities.

314 **Method 1: Distance method**

315 Figure 5A shows the weights (those > 5) for the exponential distance weight matrix. Since the entries
316 are calculated based only on the Euclidean distance between the district municipalities (and no additional
317 information), there are no significantly large weights present. As temporal information cannot be included,
318 this method produces only a single spatial weight matrix.

319 This spatial weight matrix required the largest number of principal components, namely 14, in order to
320 explain 75% of the variation in the data. This is most likely due to the lack of any form of auxiliary data or
321 information that could be used to better describe the relationship of the district municipalities. The result of
322 hierarchical clustering on the principal component observations is given in Figure 5B.

323 **Method 2: Mobile network method**

324 Figure 6 shows the spatial weight matrix for every level of lockdown that the mobile phone data spans at
administrative level 3. This spatial weight matrix identifies very strong spatial associations over relatively

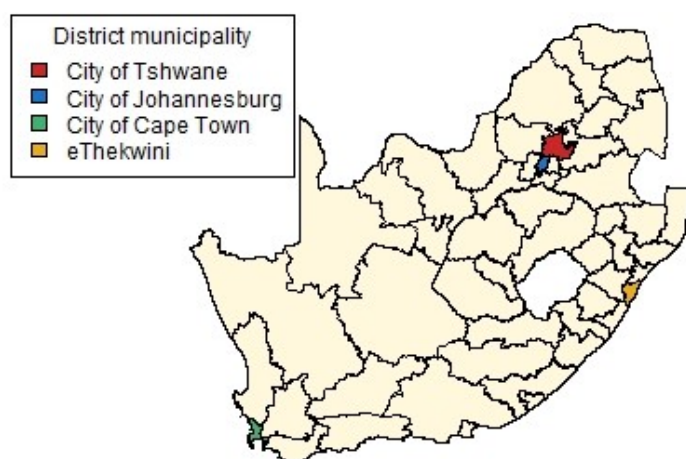


Figure 4. South Africa's district municipality boundaries and locations of four largest cities

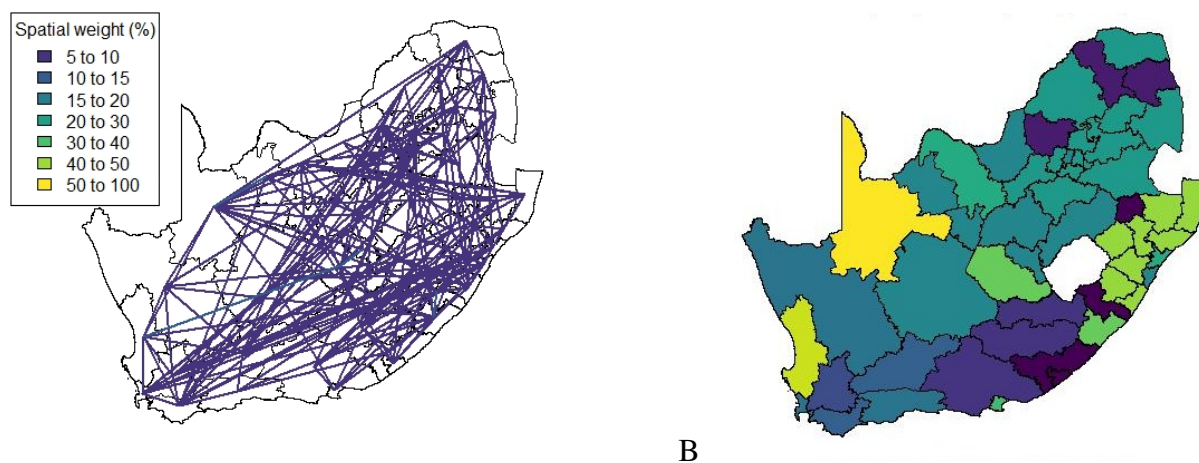


Figure 5. Method 1 A. Spatial weights (weights ≤ 5 not shown), B. Complete linkage clustering (14 clusters indicated by colours)

325 shorter distances (indicated by the yellow lines). These strong correlations appear to cluster around the
 326 edges of the country, with locations in the centre of the country displaying less spatial association overall.

327 We note that there are strong spatial associations that do not appear to be associated with any of the four
 328 major cities in the country. In particular, we note strong associations in the North-Western region of the
 329 country as well as some spatial associations across Lesotho (a neighbouring country that is landlocked by
 330 South Africa, shown in Figure 6D).

331 The spatial weight matrices for the mobile network data were also aggregated to administrative level
 332 2, shown at Figure 7. While some strong spatial associations can still be identified around the country's
 333 borders, many previously identified associations (including several significant associations spanning across
 334 the neighboring country of Lesotho) are now negligible. It is clear that while this lower spatial resolution
 335 does capture some of the spatial associations present in the data, much information is lost when aggregating
 336 between spatial resolutions.

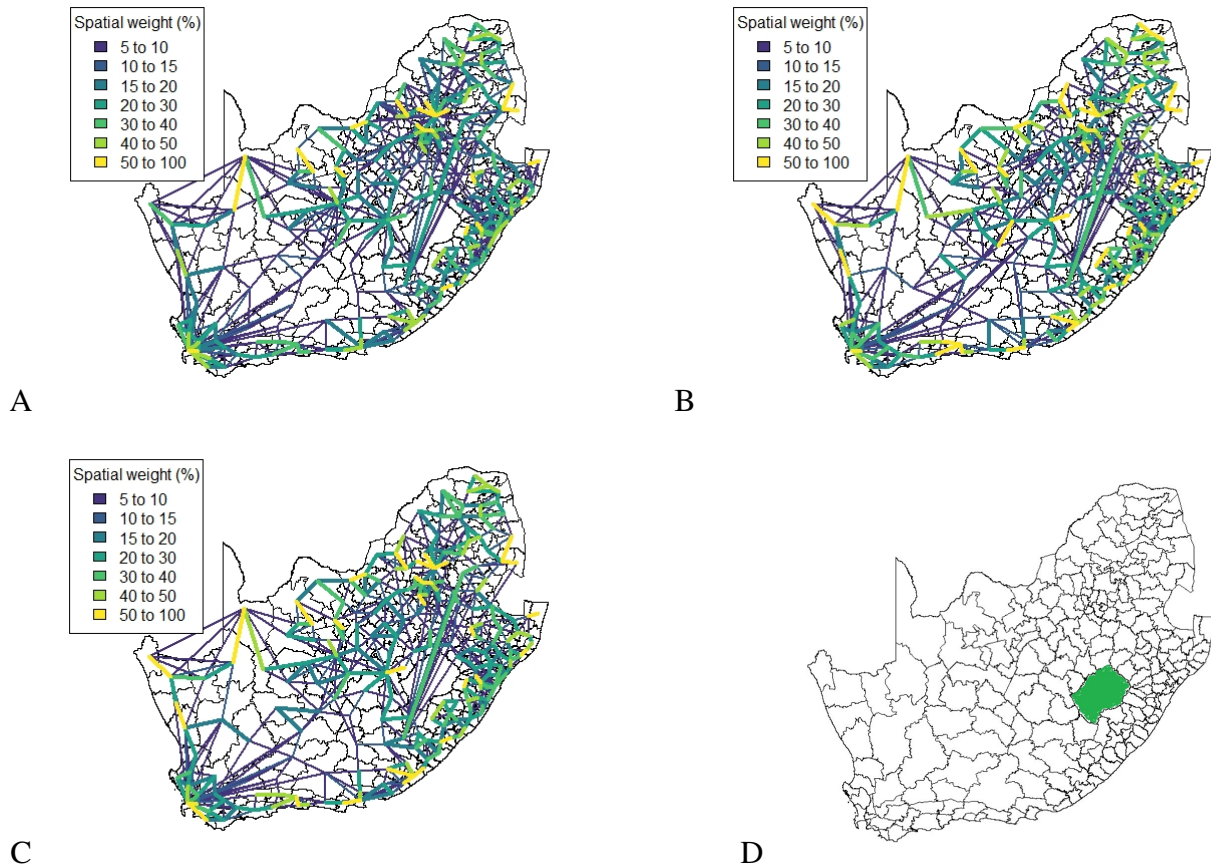


Figure 6. Method 2 spatial weight matrix entries (weights ≤ 5 not shown) A. Business as usual, B. Level 5, C. Level 4, and D. South Africa at Administrative level 3 (neighboring country Lesotho in green)

337 A notable drawback of data being at such a high spatial resolution is that it becomes very difficult to
 338 cluster locations in a meaningful way. At administrative level 3 there are 213 spatial units to consider. In
 339 order to explain just 75% of the variation in this data one requires approximately 70 principal components.
 340 Such a high number of clusters does not lend itself to easy interpretation and thus it is necessary to aggregate
 341 to a lower spatial resolution to render analysis feasible. When aggregating to administrative level 2 we
 342 find that 20 principal components are required to retain 75% of the variation present in the data. This is
 343 most likely due to the fact that the mobile network exhibits far greater daily variation than our data sources.
 344 Figure 8 shows the clustering solution.

345 **Method 3: Weighted Facebook data method**

346 This matrix construction incorporates both the Facebook population mobility data and the population
 347 size for each district municipality into the spatial weights for each district municipality pair. Figure 9
 348 shows the resulting matrix for each level of lockdown. By allowing both mobility and population size to
 349 play a role in this matrix, the strong spatial association between the four largest cities in South Africa
 350 is identified, despite the large geographical distance between them. If only Euclidean distance had been
 351 taken into account, this association would have been missed, as with Method 1. This spatial weight matrix
 352 required 9 principal components to explain 75% of the variation in the data. Figure 10 shows the results of
 353 applying hierarchical clustering to the principal component observations.

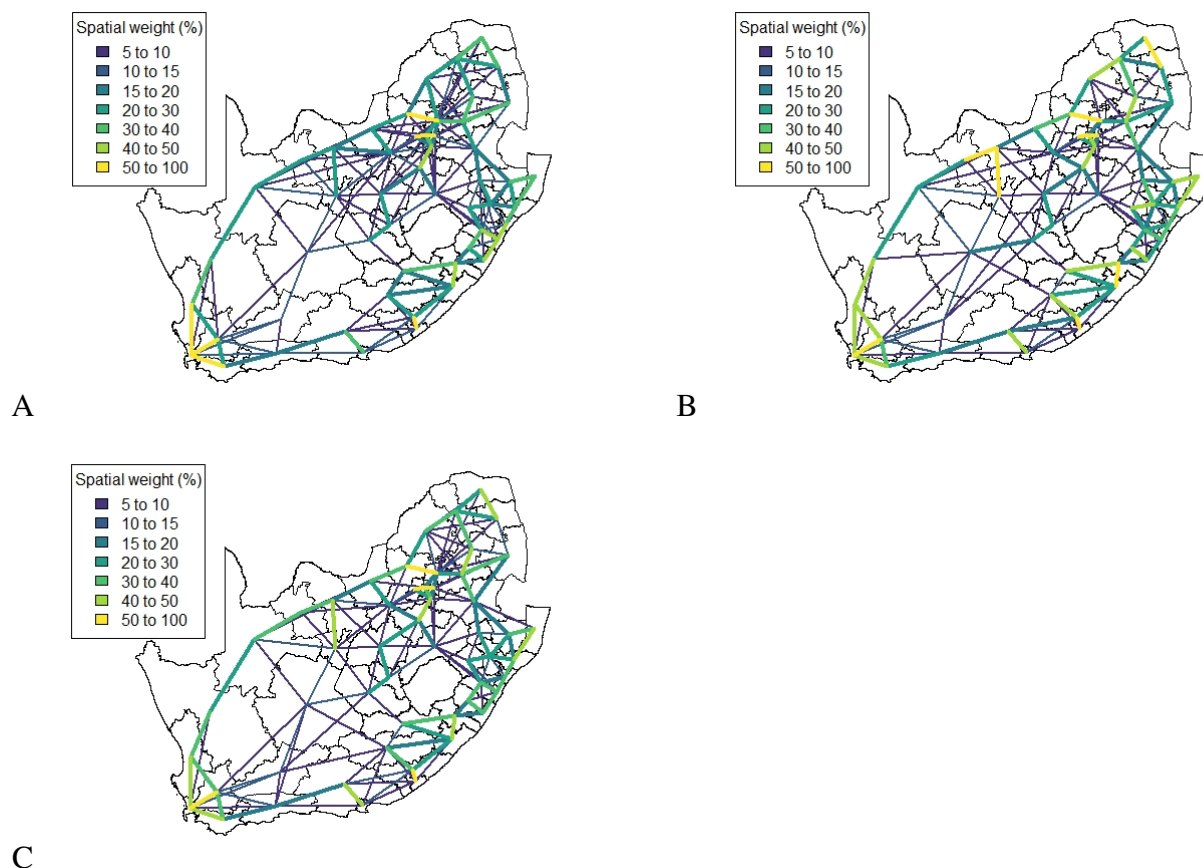


Figure 7. Method 2 spatial weight matrix entries (weights ≤ 5 not shown) A. Business as usual, B. Level 5, C. Level 4 at district municipality level

354 **Method 4: Scaled Facebook data method**

355 This spatial weight matrix was constructed as a potentially more realistic alternative to the exponential
 356 distance matrix. Despite containing a temporal element (in the form of daily mobility measurements
 357 retrieved from the Facebook data), the results for this matrix do not show any significant change across
 358 the various levels of lockdown. Figure 11 visualises the spatial weight matrix. Clustering performed on
 359 this matrix was more successful and intuitive. Only 7 components were required to explain 75% of the
 360 variation in the data. Figure 11 shows the clustering solution.

5 DISCUSSION

361 The results in Section 4 illustrate a number of ways to construct spatial weight matrices from mobility
 362 data. For the standard exponential distance method (Method 1), it is clear from Figure 5 that the clustering
 363 solution on this spatial weight matrix is not ideal. There are far too many clusters and the clustering solution
 364 reveals no clear interpretation. Although the initial matrix construction used only the distances between
 365 district municipalities, district municipalities that were located closer together were not generally clustered
 366 together.

367 The entries of the spatial weight matrix constructed using the mobile network data (Method 2), shown in
 368 Figures 6 and 7, reveal strong spatial associations over relatively short distances. The four focal largest
 369 cities in the country are clearly identified as hubs for high mobility but there are other regions, particularly
 370 those situated on or near the borders of the country, that showcase highly concentrated mobility. A possible

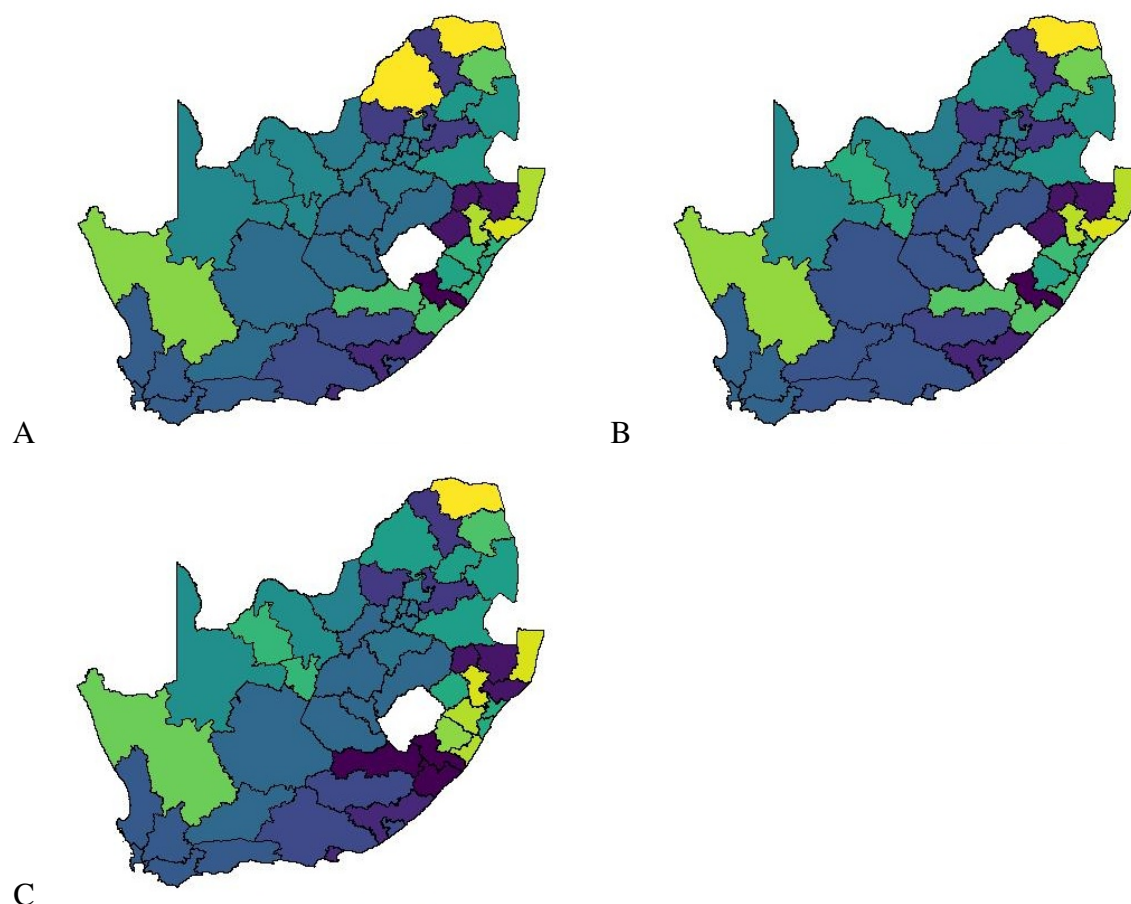


Figure 8. Method 2 complete linkage clustering results (20 clusters) A. Business as usual, B. Level 5 and C. Level 4

371 explanation for these strong spatial associations being observed far away from cities is the existence of
 372 mining activity in these areas. Given that South Africa has a very large and widespread mining sector, it
 373 seems only reasonable that any model with a spatial element should strive to incorporate these associations.
 374 The clustering solution for this spatial weight matrix, shown in Figure 8, is distinct from the other solutions
 375 in this paper in that distance is clearly not a key role player in deciding which spatial units are clustered
 376 together. Many spatial units that are situated close to one another in geographical space are not clustered
 377 together, and some spatial units are even placed into their own clusters despite having many spatial
 378 neighbours. It can be argued that this clustering solution is a more realistic reflection of the amount of
 379 travel between spatial units. The reason for this is that locations being situated closer together does not
 380 always imply that there is a higher degree of travel between these locations. The strong local connectivities
 381 picked up by this method are useful for epidemiological modelling, for example, prediction of case number
 382 hotspot movement into spatial units of higher likelihood of mobility.

383 The four largest cities in South Africa are Tshwane, Johannesburg, Cape Town and Durban, situated
 384 in the City of Tshwane, City of Johannesburg, eThekweni and City of Cape Town district municipalities
 385 respectively, as shown in Figure 4. The results in Figure 9 (method 3) show a large spatial association
 386 between these locations prior to the implementation of level 5 lockdown. Under level 5 restrictions, when
 387 the spatial influence of most district municipalities decreased, the spatial influence between these four
 388 locations became more pronounced by comparison. This most likely indicates that while smaller district

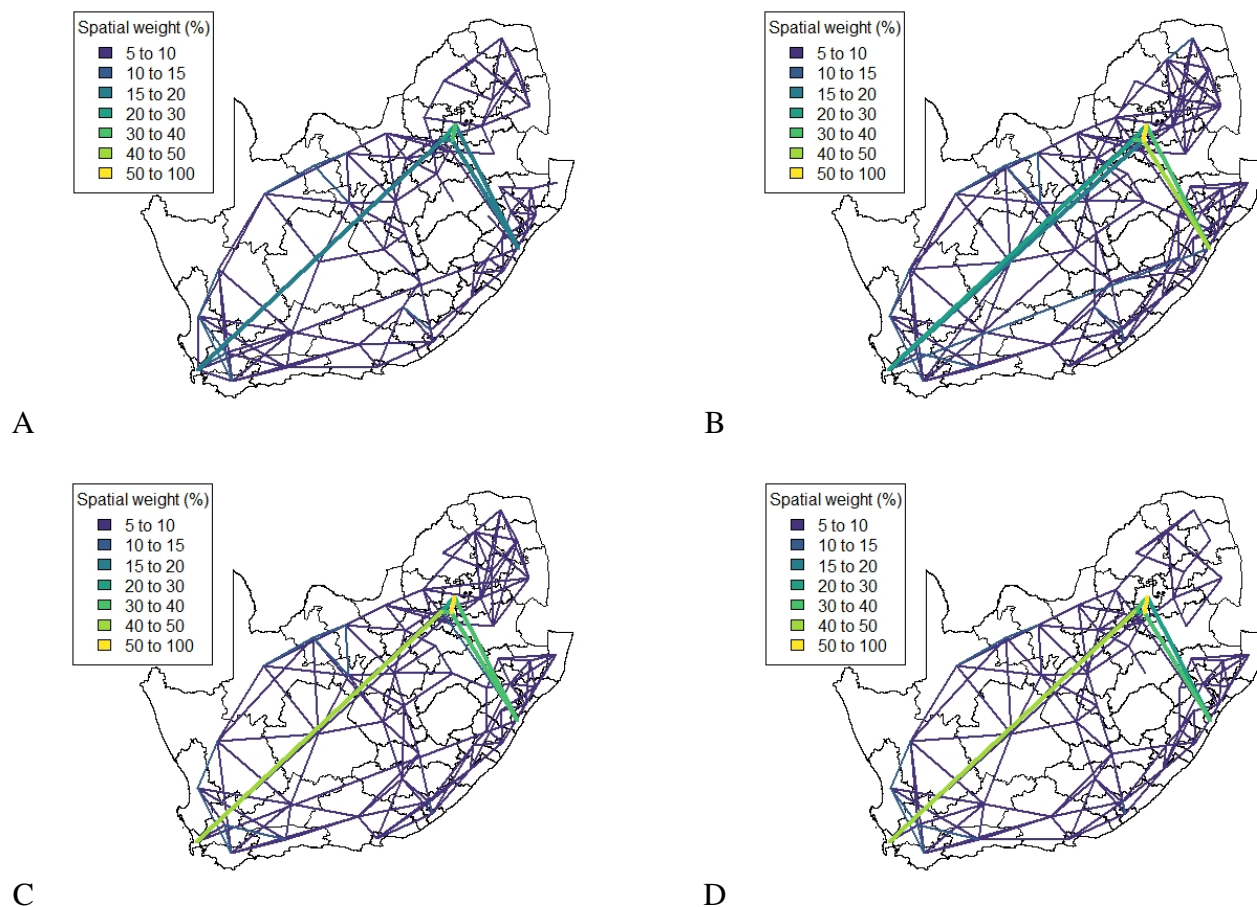


Figure 9. Method 3 spatial weight matrix entries (weights ≤ 5 not shown) A. Business as usual, B. Level 5, C. Level 4, and D. Level 3

389 municipalities were less active due to restrictions, these four were comparatively more active and still saw
 390 a sizable amount of travel between them. This seems feasible, given that these locations are the focal points
 391 for economic activity in the country and thus could not reasonably become “immobile”. As restrictions
 392 were lifted, these spatial weights were still significantly larger than those for other district municipalities,
 393 indicating that, despite restrictions being eased, the spatial influence between these four places is still
 394 significantly stronger than before the lockdown. It is also apparent that the spatial influence between less
 395 influential district municipalities has not returned to the level that they were during business as usual
 396 (pre-lockdown). Figure 10 shows that the district municipalities housing the four largest cities are all either
 397 clustered together or in clusters of their own. Other district municipalities are generally clustered together
 398 based on the distance between them. This clustering solution indicates that the four largest cities are
 399 significantly different from the locations around them. This spatial weight matrix is thus able to pinpoint
 400 the fact that these locations play a potentially larger role in spatially-dependent phenomena such as the
 401 spread of a virus. The effect in epidemiological modelling allows for longer range spatial dependency, for
 402 example, spread of the virus by daily flights between major city hubs. This is not captured by Method 2.

403 The clustering results for Method 4, shown in Figure 11, do not display any significant changes over the
 404 various levels of lockdown. Figure 11 also shows that the clusters that are formed for this spatial weight
 405 matrix are clearly based primarily on distance, but illustrates that the auxiliary Facebook data aids in
 406 constructing more finite and sensible clusters. Interestingly, we notice a district municipality that has been

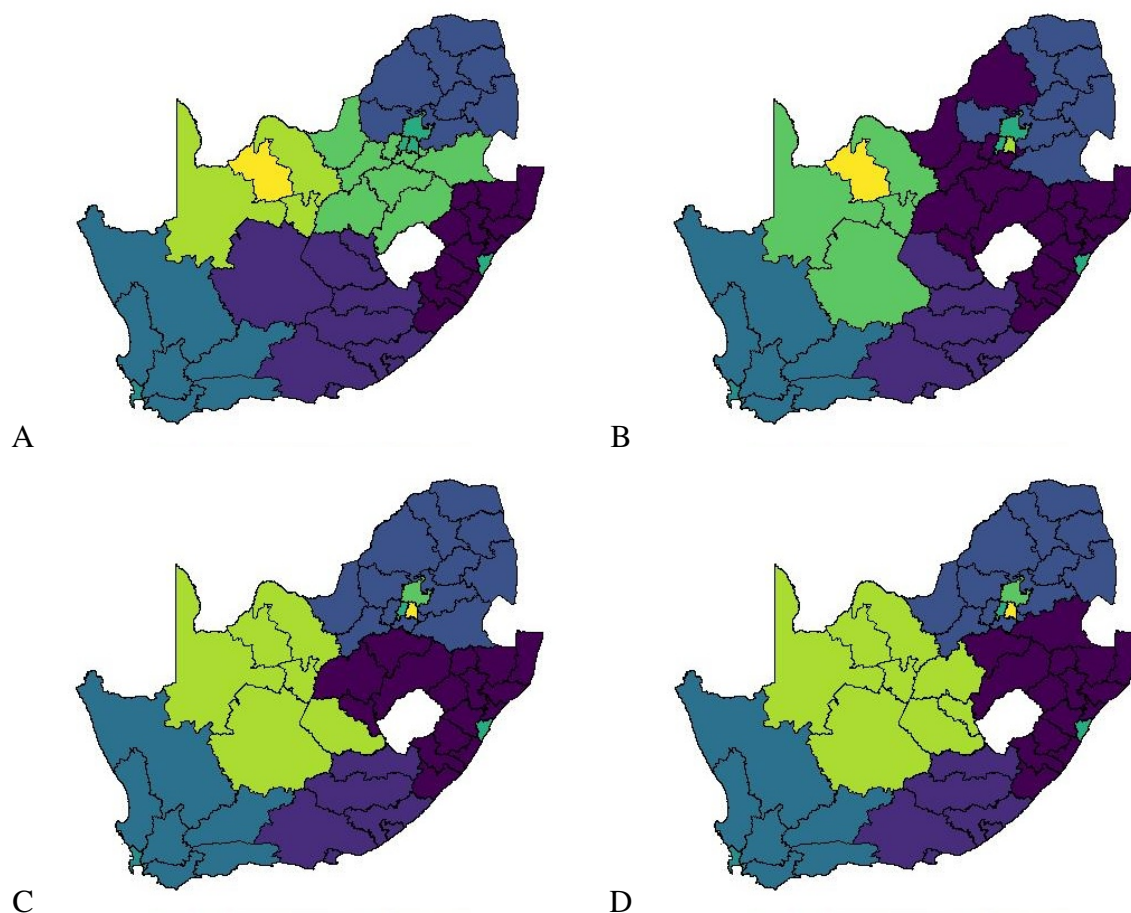


Figure 10. Method 3 complete linkage clustering results (9 clusters) A. Business as usual, B. Level 5, C. Level 4, and D. Level 3

407 classified into a cluster on its own. When inspecting the results for the other spatial weight matrices we note
 408 that this district municipality has previously also been identified as its own cluster and was shown to have
 409 strong spatial associations for Method 2. Upon further inspection we note this district municipality houses
 410 several mines. Similarly to Method 2, this spatial weight matrix is able to identify location associations that
 411 go unnoticed when relying on simple concepts such as Euclidean distance. This method may not be useful
 412 alone in epidemiological modelling and should most likely be used in conjunction with either Method 2 or
 413 3.

414 This paper shows that different representations of spatial data can offer a variety of insights and capture
 415 different relationships in the data. For example, the spatial weight matrix created using Method 3 data
 416 emphasises the prominent role of focal points in population activity. However, the spatial weight matrix
 417 constructed using Method 4 offers a scaled and smoothed way to use distance to indicate which locations
 418 have a higher spatial influence on one another. These two spatial weight matrices use the same spatial
 419 data (i.e. the Facebook for good data), but offer vastly different interpretations of spatial influence. Finally,
 420 the interpretations that were able to be made from the mobile phone data indicates that there are many
 421 potentially strong spatial associations at shorter distances that can only be identified when inspecting data
 422 at a high spatial resolution. Table 4 provides a summary of the methods used in this paper, their strengths
 423 and weaknesses, and their usability based on the results.

424 Each of these representations can be seen as valid and are complementary with regards to the insight they

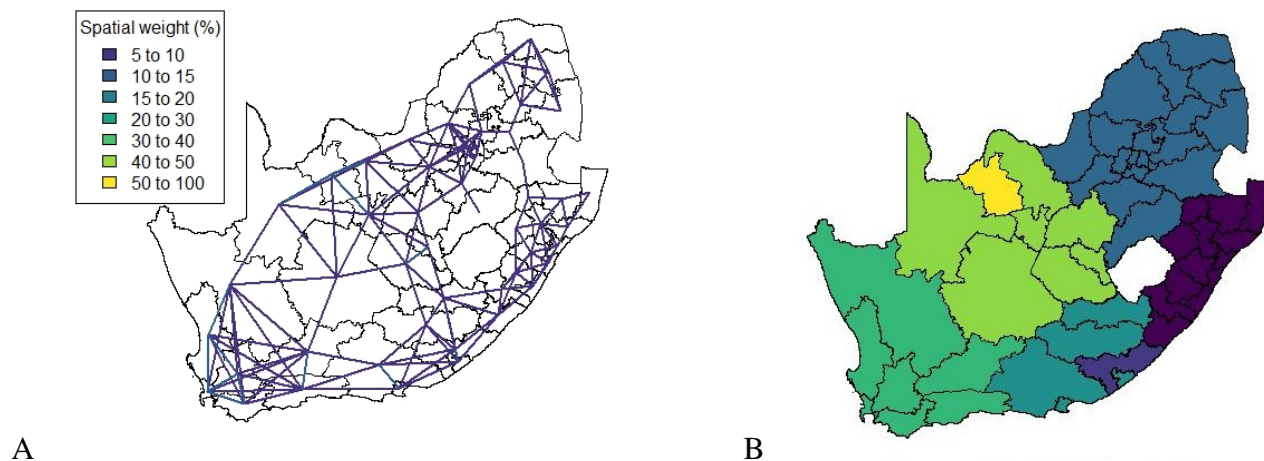


Figure 11. Method 4 A. Spatial weights (weights ≤ 5 not shown), B. Complete linkage clustering (7 clusters indicated by colours)

425 offer. Depending on the specific phenomenon under study, an argument could be made their usability based
 426 on observed patterns from the results, as in the case of a pandemic such as COVID-19, which affects not
 427 only congregated communities but allows for consequences to be felt across an entire country.

Table 4. Spatial weight matrices comparison

<i>Spatial weight matrix</i>	<i>Pro</i>	<i>Con</i>	<i>Interpretation/Contribution</i>
Method 1 - Distance	Simple to construct and understand Used often in literature	Less realistic Inadequate for clustering Lacks temporal element	Convenient to use and easy to understand and interpret. Not realistic enough for real insight.
Method 2 - Mobile network	High spatial resolution Large amounts generated passively by mobile device users	Computationally expensive Difficult to obtain Not representative Privacy concerns	Captures strong spatial associations over relatively short distances. Allows for the identification of patterns potentially missed by other methods.
Method 3 - Weighted Facebook data	Freely available data Potentially more representative	Low spatial resolution Lacks specificity	Captures association between focal points of human activity regardless of distance.
Method 4 - Scaled Facebook data	Simple to construct and understand Freely available data Potentially more representative	Lacks temporal elements Low spatial resolution	Adds additional information to previously simplistic model. Additional information improves clustering.

428 Understanding mobility during the current pandemic is essential. Both the reduction in mobility as
 429 well as retained mobility need to be well understood, and depend on reliable data collection. As shown
 430 here, data are collected in different ways and are also made available in a variety of formats. Mobility
 431 is distributionally different across strata of a region's demographics, with more mobile locations likely

432 to result in higher disease transmission. Higher resolution mobility data is important to capture these
433 differences in more detail. Even so, the spatial resolution at district municipality captures these nuances of
434 the movement under each lockdown level, and shows that significant movement still took place due to the
435 vulnerability of a large portion of South Africa's population.

436 The possibility of micro-spatial estimation (small area estimation) is something to investigate further.
437 Making use of demographic covariates, transport networks and as well as mobile network coverage maps
438 could provide connectivity matrices at higher spatial resolution, ideally at ward level. Estimation at higher
439 spatial resolution could be done by making use of a number of lower spatial resolution sources. This
440 allows for micro-scale modelling of COVID-19 spread and will allow for privacy while increasing spatial
441 resolution and providing deeper coverage in a region. Google mobility data is also available⁸ but only at
442 provincial level (administration level 1) for South Africa. This spatial resolution is too low to consider
443 estimation down to ward level, especially if alternative mobility data is available at administrative level 2.
444 However, one could also combine mobility data at different spatial resolutions in a way that takes advantage
445 of the strengths of each dataset.

446 The computational aspects of dealing with mobility data should not be overlooked. Spatial weight
447 matrices can become very large, depending on the number of spatial regions under consideration. Herein
448 the matrices were not sparse, meaning that sparse representations could not be used. Sparse representations
449 could be investigated for high spatial resolution modelling.

450 To quantify the similarity between the different spatial weight matrices, one might consider the use
451 of simple parametric measures of correlation such as Pearson's correlation coefficient. However, given
452 that there are a total of 52 spatial units (at a district municipality level) and the weights between many
453 spatial unit pairs are negligible, the spatial weight matrices can be regarded as zero-inflated. In addition to
454 making no allowance for the spatial nature of the data, namely the spatial dependency, standard measures
455 of correlation would also deliver biased results. Future research could investigate methods for comparison
456 of spatial weight matrices via appropriate correlation calculations or other techniques.

6 CONCLUSION

457 COVID-19 spreads spatially and thus the importance of mobility data for COVID-19 modeling should not
458 be disregarded. Ideally, the raw data from the mobile network providers and Facebook, if available, could
459 provide individual movements, allowing for accurate construction of spatial weight matrices. This data
460 could be anonymised and shared. However, instead the methods proposed here can be made use of. The
461 use of movement data in epidemiology is becoming an important covariate to include, without which the
462 spread can only be modelled in isolated regions. Social interactions between human beings are unavoidable.
463 Simple spatial weight matrix construction techniques, such as only taking into account distances, are not
464 always ideal when the spatial associations being captured are dependent on covariates which are not only
465 proximity based. This is made clear by the observed poor performance of Method 1 when it was used as
466 the basis of clustering. The methods presented herein and the results shown also enable epidemiological
467 modellers in considering how to incorporate spatial relationships in models. This is seldom done due
468 to limited mobility information as well as modelling complexities it introduces. However, the improved
469 accuracy in model outcomes will ultimately balance out computational complexities. The paper provides
470 insights into mobility data availability, representability as well as construction for use in spatial modelling.
471 Future research should investigate estimation to a higher spatial resolution using multiple data sources as
472 well as the effect of spatial resolution in spatial epidemiological modelling.

⁸ <https://www.google.com/covid19/mobility/> (Accessed May 2021)

ACKNOWLEDGEMENTS

473 We thank Gerbrand Mans from the CSIR for data aggregation assistance as well a Bruce Medallo from
474 WITS for mobility data advice.

475 The financial assistance of the National Research Foundation (NRF) towards this research is hereby
476 acknowledged. Opinions expressed and conclusions arrived at, are those of the authors and are not
477 necessarily to be attributed to the NRF.

478 This research is also funded by Canada's International Development Research Centre (IDRC) (Grant No.
479 109559-001).

480 The use of the Centre for High Performance Computing (www.chpc.ac.za) also made this work
481 possible.

DATA AVAILABILITY STATEMENT

482 The mobile network data used in this study is not directly available without approval, so cannot be shared
483 directly with the paper. Thank you to the NICD, South Africa for providing access this data for the
484 COVID-19 response in South Africa.

AUTHOR CONTRIBUTION

485 Conceptualisation: All; Data Curation: AP, IFR, ZK, PD; Formal Analysis: AP, IFR; Funding Acquisition:
486 PD; Investigation: AP, IFR; Methodology: AP, IFR; Project Administration: ZK, PD; Writing – original
487 draft: AP, IFR; Writing – review editing: All.

CONFLICTS OF INTEREST

488 Author Sibusiswe Khuluse Makhanya is employed by IBM, South Africa. The remaining authors declare
489 that the research was conducted in the absence of any commercial or financial relationships that could be
490 construed as a potential conflict of interest.

APPENDIX

491 Facebook for good data calculation

492 Let u represent a single individual and $U_{t,i}$ represent district municipality i at time t . The total number of
493 Bing tiles visited by inhabitants of district municipality i is then

$$\text{total_tiles}(U_{t,i}) = \sum_{u \in U_{t,i}} \min(\text{tiles}(u), 200).$$

494 Note that the maximum number of Bing tiles visited that a single individual can contribute is restricted to
495 200. In order to preserve user privacy, an error term was included by drawing from a Laplace distribution
496 with parameters 0 and $\frac{F}{\epsilon}$ where F = sensitivity parameter and ϵ = noise parameter as follows

$$\text{total_tiles}'(U_{t,i}) = \text{total_tiles}(U_{t,i}) + \text{Laplace}\left(0, \frac{F}{\epsilon}\right).$$

497 The average number of tiles per district municipality was then calculated as

$$\text{avg_tiles}'(U_{t,i}) = \frac{\text{total_tiles}'(U_{t,i})}{|U_{t,i}|}.$$

498 The mobility value for each district municipality and for each day was then finally expressed with respect
499 to the baseline as

$$F_i^{(t)} = \frac{\text{avg_tiles}(U_{t,i}) - \text{baseline_avg_tiles}'(i, \text{day_of_the_week}(t))}{\text{baseline_avg_tiles}'(i, \text{day_of_the_week}(t))}.$$

500 For further details regarding this data see [https://research.fb.com/blog/2020/06/protecting-privacy-in-](https://research.fb.com/blog/2020/06/protecting-privacy-in-facebook-mobility-data-during-the-covid-19-response/)
501 [facebook-mobility-data-during-the-covid-19-response/](https://research.fb.com/blog/2020/06/protecting-privacy-in-facebook-mobility-data-during-the-covid-19-response/).